

Variance V.S. Bias

Mathematically:

Let the variable we are trying to predict as y and other covariates as x . We assume there is a relationship between the two such that

$$y = f(x) + e$$

where e is the error term and it's normally distributed with a mean of 0.

We will make a model $\hat{f}(x)$ of $f(x)$ using linear regression or any other modeling technique.

So, the expected squared error at a point x is

$$\text{Err}(x) = E[(y - \hat{f}(x))^2]$$

The $\text{Err}(x)$ can be further decomposed as

$$\text{Err}(x) = (E[\hat{f}(x)] - f(x))^2 + E[(\hat{f}(x) - E[\hat{f}(x)])^2] + \sigma_e^2,$$

$$\text{Err}(x) = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}.$$

Irreducible error is the error that can't be reduced by creating good models. It is a measure of the amount of noise in our data.

Here it is important to understand that no matter how good we make our model, our

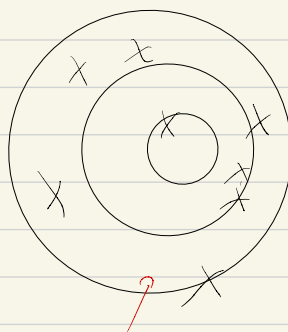
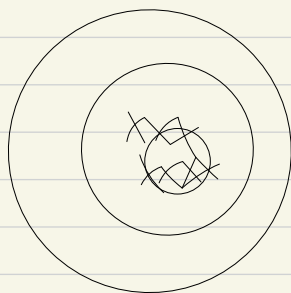
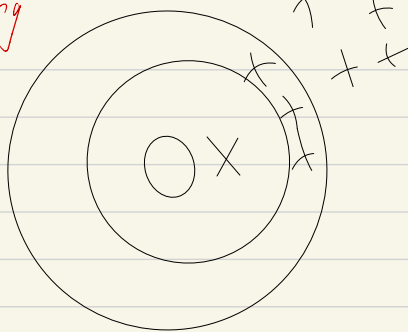
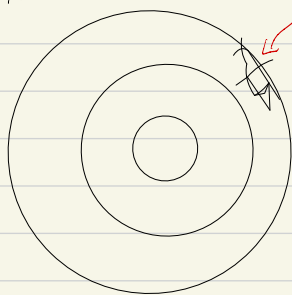
Low
Variance

underfitting

high variance

High
Bias

Low
Bias



overfitting

Why is Bias Variance Tradeoff?

If our model is too simple and has very few parameters then it may have high bias and low variance. On the other hand if our model has large number of parameters then it's going to have high variance and low bias. So we need to find the right/good balance without overfitting and underfitting the data.

This tradeoff in complexity is why there is a tradeoff between bias and variance. An algorithm